
O²TD: (Near)-Optimal Off-Policy TD Learning

Bo Liu
Auburn University
boliu@auburn.edu

Daoming Lyu
Auburn University
dzt0053@auburn.edu

Wen Dong
University of Buffalo
wendong@buffalo.edu

Saad Biaz
Auburn University
biazsaa@auburn.edu

Abstract

Temporal difference learning and Residual Gradient methods are the most widely used temporal difference based learning algorithms; however, it has been shown that none of their objective functions is optimal w.r.t approximating the true value function V . Two novel algorithms are proposed to approximate the true value function V . This paper makes the following contributions:

- A batch algorithm that can help find the approximate optimal off-policy prediction of the true value function V .
- A linear computational cost (per step) near-optimal algorithm that can learn from a collection of off-policy samples.
- A new perspective of the emphatic temporal difference learning which bridges the gap between off-policy optimality and off-policy stability.

1 Introduction

Temporal difference (TD) learning is a widely used method in reinforcement learning. There are two fundamental problems in temporal difference learning. The **first problem** is the off-policy stability. Although TD converges when samples are drawn “on-policy” by sampling from the Markov chain underlying a policy in a Markov decision process, it can be shown to be divergent when samples are drawn “off-policy.” Off-policy stable methods are of wider applications since they can learn while executing an exploratory policy, learn from demonstrations, and learn multiple tasks in parallel. The **second problem** is the optimality with function approximation. An accurate prediction of the value function will greatly help improve the policy optimization, which is the ultimate

goal of reinforcement learning tasks. On the other hand, a bad value function prediction will lead to a low-quality policy [Sutton and Barto, 1998].

Several different approaches have been explored to address the problem of off-policy temporal difference learning. Baird’s residual gradient (**RG**) method [Baird, 1995] is the first approach with linear complexity per step, but it requires double sampling and also converges to an inferior solution. Gordon [1996] proposed the “averager” method, which needs to store many training examples, and thus is not practical for large-scale applications. The off-policy LSTD [Yu, 2010] is off-policy convergent, but its per-step computational complexity is quadratic in the number of parameters d of the function approximator. Sutton *et al.* [2008, 2009] proposed the family of gradient-based temporal difference (**GTD**) algorithms which are proven to be asymptotically off-policy convergent using stochastic approximation [Borkar, 2008].

Another direction of temporal difference learning, optimal temporal difference learning, seems to draw relatively insufficient attention. It is well-known that the asymptotic solutions of TD and GTD are not the true value function V , but the solution of a projected fixed point equation [Sutton *et al.*, 2009]. On the other hand, the residual gradient method converges to another solution, which is often inferior to the TD solution. However, as pointed out by Scherrer [2010], both the TD and residual gradient method can be unified as the oblique projection of the true value function V with different oblique projection directions, and *neither* of them is optimal in the sense of approximating the true value function V . To the best of our knowledge, the most relevant to our work is the optimal Dantzig Selector TD learning [Liu *et al.*, 2016], which aims to find the best denoising matrix for the purpose of feature selection, when the number of samples n is much larger than the number of function approximation parameters d .

This paper attempts to improve the prediction of value

function based on the technique of oblique projection. Here is a roadmap for the rest of the paper. Section 3 introduces the relationship between the optimal approximation of the true value function V with the oblique projected fixed point equations, which reduces the problem to finding the optimal oblique projection direction. Unfortunately, this cannot be directly computed. To this end, Section 4 proposes an approximation criterion and two algorithms, i.e., a state-aggregated batch algorithm and a state-weighted stochastic algorithm. Related work is discussed in Section 5. Section 6 presents the experimental results evaluating the effectiveness of the proposed approaches.

2 Preliminary

Reinforcement Learning (**RL**) [Bertsekas and Tsitsiklis, 1996; Sutton and Barto, 1998] is a class of learning problems in which an agent interacts with an unfamiliar, dynamic, and stochastic environment, where the agent’s goal is to optimize some measure of its long-term performance. This interaction is conventionally modeled as a Markov decision process (**MDP**). An MDP is defined as the tuple $(\mathcal{S}, \mathcal{A}, P_{ss'}^a, R, \gamma)$, where \mathcal{S} and \mathcal{A} are the sets of states and actions, the transition kernel $P_{ss'}^a$ specifying the probability of transition from state $s \in \mathcal{S}$ to state $s' \in \mathcal{S}$ by taking action $a \in \mathcal{A}$, $R(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function bounded by R_{\max} , and $0 \leq \gamma < 1$ is a discount factor. A stationary policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a probabilistic mapping from states to actions. The main objective of a RL algorithm is to find an optimal policy. In order to achieve this goal, a key step in many algorithms is to calculate the value function of a given policy π , i.e., $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$, a process known as *policy evaluation*. It is known that V^π is the unique fixed-point of the *Bellman operator* T^π , i.e.,

$$V^\pi = T^\pi V^\pi = R^\pi + \gamma P^\pi V^\pi, \quad (1)$$

where R^π and P^π are respectively the reward function and transition kernel of the Markov chain induced by policy π . In Eq. 1, we may think of V^π as a $|\mathcal{S}|$ -dimensional vector and write everything in vector/matrix form. We also denote $L^\pi := I - \gamma P^\pi$. In the following, to simplify the notation, we often drop the dependence of T^π , V^π , R^π , and P^π to π .

We denote by π_b , the behavior policy that generates the data, and by π , the target policy that we would like to evaluate. They are the same in the on-policy setting and different in the off-policy scenario. For each state-action pair (s_i, a_i) , such that $\pi_b(a_i|s_i) > 0$, we define the importance-weighting factor $\rho_i = \pi(a_i|s_i)/\pi_b(a_i|s_i)$ with $\rho_{\max} \geq 0$ being its maximum value over the state-action pairs.

When \mathcal{S} is large or infinite, we often use a linear approximation architecture for V^π with parameters $\theta \in \mathbb{R}^d$ and K -bounded basis functions $\{\varphi_i\}_{i=1}^d$, i.e., $\varphi_i : \mathcal{S} \rightarrow \mathbb{R}$ and $\max_i \|\varphi_i\|_\infty \leq K$. We denote by $\phi(\cdot) := (\varphi_1(\cdot), \dots, \varphi_d(\cdot))^\top$ the feature vector and by \mathcal{F} the linear function space spanned by the basis functions $\{\varphi_i\}_{i=1}^d$, i.e., $\mathcal{F} = \{f_\theta \mid \theta \in \mathbb{R}^d \text{ and } f_\theta(\cdot) = \phi(\cdot)^\top \theta\}$. We may write the approximation of V in \mathcal{F} in the vector form as $\hat{v} = \Phi\theta$, where Φ is the $|\mathcal{S}| \times d$ feature matrix, and we denote

$$\Delta := (I - \gamma P^\pi)\Phi = L^\pi\Phi. \quad (2)$$

When only n training samples of the form $\mathcal{D} = \{(s_i, a_i, r_i = r(s_i, a_i), s'_i)\}_{i=1}^n$, $s_i \sim \xi$, $a_i \sim \pi_b(\cdot|s_i)$, $s'_i \sim P(\cdot|s_i, a_i)$, are available (ξ is a vector representing the probability distribution over the state space \mathcal{S}), we denote by $\delta_i(\theta) := r_i + \gamma\phi'_i{}^\top\theta - \phi_i^\top\theta$, the TD error for the i -th sample (s_i, r_i, s'_i) and define $\Delta\phi_i = \rho_i(\phi_i - \gamma\phi'_i)$. We also denote the sample-based state-aggregated estimation of Δ (resp. R), termed as $\hat{\Delta}$ (resp. \hat{R}), i.e., given sample set \mathcal{D} , the i -th and j -th samples are aggregated if $s_i = s_j$, which is a standard approach used in state aggregation methods [Singh *et al.*, 1995]. Finally, we define the matrices C as $C := \mathbb{E}[\phi_i\phi_i^\top]$, where the expectations are w.r.t. ξ and P^{π_b} . We also denote by Ξ , the diagonal matrix whose elements are $\xi(s)$, and $\xi_{\max} := \max_s \xi(s)$. For each sample i in the training set \mathcal{D} , the unbiased estimate of C is $\hat{C}_i := \phi_i\phi_i^\top$.

3 Problem Formulation

This section presents the motivation of this research, i.e., exploring the possible optimal value function approximation in a model-free reinforcement learning setting.

It is evident that given the functional space \mathcal{F} and the approximation of V in \mathcal{F} in the vector form represented as $\hat{v} = \Phi\theta$, the optimal approximation is $v^* = \Pi V$, where $\Pi = \Phi(\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi$ is the weighted least-squares projection weighted by the state distribution Ξ . This is obtained from $\arg \min_{\hat{v}} \|\hat{v} - V\|_\Xi^2$. It is also well-known that the TD solution \hat{v}_{TD} does not converge to v^* but to the unique fixed-point solution of $\hat{v} = \Pi T \hat{v}$. Several intuitive questions arise here, such as

1. What is the approximation error bound between \hat{v}_{TD} and V ?
2. What is the relation of representation between \hat{v}_{TD} and V , i.e., if \hat{v}_{TD} can be analytically represented by V ?

The first question has been answered in [Tsitsiklis and Van Roy, 1997], where an upper bound was given as $\|V - \hat{v}_{TD}\|_\xi \leq \frac{1}{\sqrt{1-\gamma^2}}\|V - v^*\|_\xi$. The answer to the second question requires the notion of oblique projection defined in Section 3.1.

3.1 Oblique Projection and Optimal Projection

The oblique projection tuple (Φ, X) is defined as follows, where X is a matrix with the same size as Φ .

Definition. The *Oblique Projection* operator Π_Φ^X is defined as

$$\Pi_\Phi^X = \Phi(X^\top \Phi)^{-1} X^\top, \quad (3)$$

which specifies a projection orthogonal to $\text{span}(X)$ and onto $\text{span}(\Phi)$. It can be easily deduced that the weighted orthogonal projection Φ can be written as $\Pi = \Pi_\Phi^{\Xi\Phi}$.

It is easy to verify that the projected fixed point equation in temporal difference learning, $\hat{v} = \Pi T^\pi(\hat{v})$, can be extended to a more general setting by extending the weighted least-squares projection operator to oblique projection operator as

$$\hat{v} = \Pi_\Phi^X T^\pi(\hat{v}), \quad (4)$$

It turns out that both TD and RG solutions are oblique projections with different X , where $X_{TD} = \Xi\Phi$, $X_{RG} = \Xi L^\pi \Phi$ [Scherrer, 2010]. One may be interested in the relation between the true value function V and the solutions of the fixed-point equation. The relation is shown in Lemma 1.

Lemma 1. [Scherrer, 2010] *The solution of the oblique projected fixed-point equation $\hat{v} = \Pi_\Phi^X T(\hat{v})$ w.r.t the oblique projection Π_Φ^X can be represented as the oblique projection $\Pi_\Phi^{L^{\pi^\top} X}$ of the true value function V , i.e.,*

$$\hat{v} = \Pi_\Phi^X T^\pi(\hat{v}) = \Pi_\Phi^{L^{\pi^\top} X} V, \quad (5)$$

where $L^\pi = (I - \gamma P^\pi)$.

Proof. Please refer to Scherrer [2010] for a detailed proof. \square

Remark: Lemma 1 helps to identify the equivalence between oblique projection of the true value function V , i.e., $\Pi_\Phi^{L^{\pi^\top} X} V$ and the solution of the oblique projected fixed-point equation, i.e., $\hat{v} = \Pi_\Phi^X T^\pi \hat{v}$. Figure 1 is an illustration of the oblique projection.

An intuitive question to ask is what the best oblique projection X is. Is it either TD, RG, or some interpolation between them, or none of the above? To answer this question, we present the following proposition, which is the workhorse of this paper.

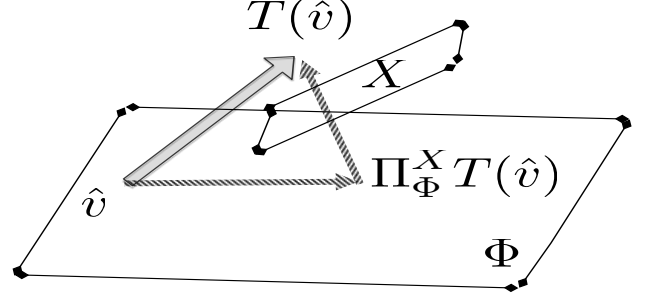


Figure 1: An Illustration of Oblique Projected TD

Lemma 2. [Scherrer, 2010] *Given Φ , if V does not lie in $\text{span}(\Phi)$, the optimal approximation is $v^* = \Pi V = \Phi(\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi V$, and the corresponding oblique projection X^* in the fixed point equation*

$$v^* = \Pi_\Phi^{X^*} T v^* \quad (6)$$

is

$$X^* = (L^{\pi^\top})^{-1} \Xi \Phi. \quad (7)$$

Proof. From Lemma 1, we know that X^* satisfies $v^* = \Pi_\Phi^{(L^{\pi^\top})^{-1} X^*} V$. Let $\Pi_\Phi^{(L^{\pi^\top})^{-1} X^*} = \Pi$, we have

$$(L^{\pi^\top}) X^* = \Xi \Phi, \quad (8)$$

and thus we can have Eq. (7), which completes the proof. \square

Although the analytical formulation of X^* is clear, it is intractable to compute. The major reason is that $(L^{\pi^\top})^{-1}$ is computational prohibitive since the exact P^π is not known. This paper will present techniques to compute X^* approximately in the following.

4 Algorithm Design

Given the knowledge of the oblique projection and the problem of the computational intractability to compute X^* , a criterion is proposed to approximate X^* . Based on this criteria, two algorithms are proposed. The first is based on state-aggregated two-stage approximation, and the second is based on state-dependent diagonalized approximation.

4.1 Approximate Criteria

Before presenting the algorithm design, we first introduce a simple but important property of the optimal projection matrix X^* . Notice since $X^* = (L^{\pi^\top})^{-1} \Xi \Phi$, and thus we have

$$\Delta^\top X^* = \Phi^\top (L^{\pi^\top}) ((L^{\pi^\top})^{-1} \Xi \Phi) = \Phi^\top \Xi \Phi = C. \quad (9)$$

Motivated by this, Proposition 1 is presented to formulate the cornerstone of this paper.

Proposition 1. *For state aggregated $\hat{\Delta}$, there is*

$$\mathbb{E}_{\pi_b}[\hat{\Delta}] = L^\pi \Phi, \quad (10)$$

and thus for the optimal oblique projection X^ and the corresponding $v^* = \Phi\theta^*$, the following holds*

$$\mathbb{E}_{\pi_b}[\hat{\Delta}]^\top X^* = \mathbb{E}_\xi[\hat{C}] \quad (11)$$

$$X^{*\top} \mathbb{E}_{\pi_b}[\hat{\Delta}]\theta^* = X^{*\top} \mathbb{E}_\xi[\hat{R}]. \quad (12)$$

Proof. Eq. (10) is derived as follows,

$$\begin{aligned} \mathbb{E}_{\pi_b}[\hat{\Delta}] &= \mathbb{E}_{\pi_b}[\rho_i \Delta \phi_i] \\ &= \sum_{a_i} \pi_b(a_i | s_i) \frac{\pi(a_i | s_i)}{\pi_b(a_i | s_i)} (\Delta \phi_i)^\top \\ &= \sum_{a_i} \pi(a_i | s_i) (\Delta \phi_i)^\top \\ &= L^\pi \Phi. \end{aligned} \quad (13)$$

Then we have

$$\begin{aligned} \mathbb{E}_{\pi_b}[\hat{\Delta}^\top] X^* &= (L^\pi \Phi)^\top (L^{\pi^\top})^{-1} \Xi \Phi \\ &= \Phi^\top \Xi \Phi = C. \end{aligned} \quad (14)$$

Insert Eq. (10), $\mathbb{E}_\xi[\hat{C}] = C$, and $\mathbb{E}_\xi[\hat{R}] = R$ into Eq. (6), there is

$$\begin{aligned} \Phi\theta^* &= \Pi_\Phi^{X^*} (R + \gamma\Phi'\theta^*) \\ &= \Phi(X^{*\top}\Phi)^{-1} X^{*\top} (R + \gamma\Phi'\theta^*) \\ \theta^* &= (X^{*\top}\Phi)^{-1} X^{*\top} (R + \gamma\Phi'\theta^*) \\ X^{*\top}\Phi\theta^* &= X^{*\top} (R + \gamma\Phi'\theta^*) \\ X^{*\top}(\Phi - \gamma\Phi')\theta^* &= X^{*\top} R \\ X^{*\top}\Delta\theta^* &= X^{*\top} R. \end{aligned}$$

This completes the proof. \square

4.2 Two-Stage State-Aggregated Batch Algorithm

Motivated by Proposition 2, the following two-stage near-optimal off-policy TD algorithm is proposed, where the first step is

$$\hat{X} = \arg \min_X \frac{1}{2} \|\hat{\Delta}^\top X - \hat{C}\|_F^2. \quad (15)$$

This problem is a well-defined convex problem, and there exists a unique solution. When $(\hat{\Delta}\hat{\Delta}^\top)$ is nonsingular, the closed-form least-squares solution is computed as

$$\hat{X} = (\hat{\Delta}\hat{\Delta}^\top)^{-1}(\hat{\Delta}C). \quad (16)$$

On the other hand, if $(\hat{\Delta}\hat{\Delta}^\top)$ is singular, which is more general, Eq. (15) can be solved via gradient descent method and can be further accelerated by Nesterov's accelerated gradient method [Nesterov, 2004]. The second step is to compute $\hat{\theta}$, i.e.,

$$\hat{\theta} = \arg \min_{\theta} \|\hat{X}^\top (\hat{\Delta}\theta - \hat{R})\|_\xi^2. \quad (17)$$

This is a well-defined convex problem, and the solution is unique and can be easily solved via gradient descent method. When $(\hat{X}^\top \hat{\Delta})$ is nonsingular, \hat{X} can be simply solved via the one-shot least-squares solution

$$\hat{\theta} = (\hat{X}^\top \hat{\Delta})^{-1} \hat{X}^\top R. \quad (18)$$

Based on these, we propose the *State-aggregated Optimal TD Algorithm (SOTD)* as follows.

Algorithm 1 State-aggregated Optimal TD Algorithm (SOTD)

- 1: INPUT: Sample set $\{\phi_i, r_i, \phi_i'\}_{i=1}^n$
 - 2: Compute $\hat{\Delta}, \hat{C}, \hat{R}$.
 - 3: Compute \hat{X} as in Eq. (15).
 - 4: Compute $\hat{\theta}$ as in Eq. (17).
-

4.3 State-Dependent Optimal Off-Policy TD learning

Algorithm 1 can find the near-optimal projection matrix, however, there is an apparent drawback of computing \hat{X} in this way because of computational complexity. Note that \hat{X} is a $|\mathcal{S}| \times d$ matrix, which is computationally costly in large-scale reinforcement learning problems where the number of states $|\mathcal{S}|$ is large, or in continuous state space. To this end, the following algorithm is designed to tackle difficulties mentioned above.

In real applications where $d \ll |\mathcal{S}|$ or the state space is continuous, the proposed algorithm would not work well in practice since it has to compute a $|\mathcal{S}| \times d$ matrix \hat{X} . A desirable way out is to approximate the $(L^{\pi^\top})^{-1}$ with a diagonal matrix Ω , such that each row of Ω does not depend on other states, but only on its corresponding state. With such an assumption, \hat{X} can be represented via a product of matrices as follows,

$$\hat{X} = \Omega \Xi \Phi, \quad (19)$$

where Ω is a $|\mathcal{S}| \times |\mathcal{S}|$ diagonal matrix. The i -th diagonal entry of Ω is denoted as ω_i , i.e., $\Omega_{ii} := \omega_i$. We term Ω as “state-dependent” diagonal matrix. Then the optimization problem reduces to

$$\hat{\Omega} = \arg \min_{\Omega} \frac{1}{2} \|\hat{\Delta}^\top \Omega \Xi \Phi - \hat{C}\|_F^2, \quad \text{s.t.} \quad \Omega_{ij} = 0, \quad i \neq j. \quad (20)$$

It is easy to prove the following

$$\hat{\Delta}^\top \Omega \Xi \Phi = \mathbb{E}[\rho_i \omega_i \phi_i \Delta \phi_i^\top]. \quad (21)$$

Based on the assumption that ω_i should be only (current) state-dependent, we have the following relaxed objective function, i.e., for the i -th sample,

$$\forall i, \omega_i = \arg \min_{\omega} \|\omega \rho_i \phi_i \Delta \phi_i^\top - \phi_i \phi_i^\top\|_F^2. \quad (22)$$

Trace norm minimization can also be used, i.e.,

$$\omega_i = \arg \min_{\omega} \|\phi_i (\omega \rho_i \Delta \phi_i - \phi_i)^\top\|_*. \quad (23)$$

Two issues arise here:

- *Computational cost.* Trace norm minimization is usually more computationally expensive since it involves the singular value decomposition (SVD) operation.
- *Choice of the norm.* The issue here is to select the best norm as the objective function. Although there is already several pieces of literature discussing this problem, however, it remains unclear that at first glance, which norm minimization would achieve the best result in our problem.

We will resolve these two concerns by scrutinizing the structure of the problem. Notice that Eq. (22) can be written as $\omega_i = \arg \min_{\omega} \|\phi_i (\omega \rho_i \Delta \phi_i - \phi_i)^\top\|_F^2$. Since $\phi_i (\omega \rho_i \Delta \phi_i - \phi_i)^\top$ is a rank-1 matrix, the solution is identical w.r.t Frobenius norm and trace norm, and the closed-form solution is

$$\omega_i = \frac{\Delta \phi_i^\top \phi_i}{\rho_i \Delta \phi_i^\top \Delta \phi_i}. \quad (24)$$

Interested readers will find a detailed deduction in the Appendix. The update law is thus as follows,

$$\theta_{i+1} = \theta_i + \alpha_i \rho_i \omega_i \delta_i \phi_i. \quad (25)$$

Samples with zero importance ratio (i.e., $\rho_i = 0$) are discarded. Now it is ready to formulate the *Optimal Off-Policy TD Algorithm* (**O²TD**) algorithm as in Algorithm 2. It is easy to verify that the computational cost per step is $O(d)$, as can be seen from the computation of Eq. (24) and (25).

5 Related Work

One of the related work to optimal temporal difference learning is the emphatic temporal difference learning (ETD) work by Sutton *et al.* [2015]. That work was motivated by the off-policy convergence issue, and we

Algorithm 2 Optimal Off-Policy TD Algorithm (O²TD)

- 1: INPUT: Sample set $\{\phi_i, r_i, \phi_i'\}_{i=1}^n$
 - 2: **for** $i = 1, \dots, n$ **do**
 - 3: Compute $\phi_i, \Delta \phi_i, \delta_i = r_i + \gamma \phi_i'^\top \theta_i - \phi_i^\top \theta_i$.
 - 4: Compute ω_i according to Eq. (24).
 - 5: Compute θ_{i+1} according to Eq. (25).
 - 6: **end for**
-

will shed new light on the algorithms from the optimality perspective. Similar to O²TD, ETD also assumes that the optimal projection X^* can be approximated by the product of diagonal matrices Ω, Ξ and the Φ matrices, i.e., the near-optimal projection matrix is formulated as in Eq. (19). Then a different technique is used based on the power series expansion, i.e.,

$$(L^\pi)^{-1} = (I - \gamma P^\pi)^{-1} = \sum_{i=0}^{\infty} (\gamma P^\pi)^i. \quad (26)$$

Then the power series expansion is used to compute $\Omega \Xi$ as a whole. Since the optimal oblique projection matrix is $X^* = (L^{\pi^\top})^{-1} \Xi \Phi$, it is evident that $\hat{X} = \Omega \Xi \Phi$ should be as close as possible to X^* , especially the diagonal elements. The diagonal elements of \hat{X} are represented as a (column) vector f . One conjecture is that for the diagonal matrix of \hat{X} , it is desired that $f = (L^{\pi^\top})^{-1} \xi$. By using the power series expansion, f can be expanded as

$$\begin{aligned} f &= (L^{\pi^\top})^{-1} \xi = \left(\sum_{i=0}^{\infty} (\gamma P^{\pi^\top})^i \right) \xi \\ &= (I + \gamma P^{\pi^\top} + (\gamma P^{\pi^\top})^2 + \dots + (\gamma P^{\pi^\top})^k + \dots) \xi. \end{aligned} \quad (27)$$

Readers familiar with the emphatic TD learning algorithm know that this is actually identical to Equation (13) in the paper by Sutton *et al.* [2015], where a scalar follow-on trace is computed as ¹

$$F_0 = 1; \quad F_t = I + \gamma \rho_{t-1} F_{t-1}, \quad t > 0, \quad (29)$$

and it turns out that

$$f_i = \xi(i) \lim_{t \rightarrow \infty} \mathbb{E}[F_t | S_t = s_i], \quad (30)$$

which will lead to the standard emphatic TD(0) algorithm,

$$\theta_{t+1} = \theta_t + \alpha_t F_t \rho_t \delta_t \phi_t. \quad (31)$$

Due to space limitations, we refer interested readers to [Sutton *et al.*, 2015] for more details of the algorithm,

¹We use subscript \bullet_t to denote sequential samples, and subscription \bullet_i to denote samples that are randomly sampled with replacement.

and [Hallak *et al.*, 2015; Yu, 2015] for more theoretical analysis. It should also be noted that although this section does not provide any further extension of the ETD algorithm regarding algorithm design and analysis, to the best of our knowledge, it is the first time associating the ETD algorithm with near optimal temporal difference learning. This sheds a helpful light in understanding the family of the emphatic TD learning algorithms and the design of the follow-on trace. However, the ETD algorithm requires sequential sampling condition, i.e., $s'_t = s_{t+1}, \forall t > 0$, which is not suitable for a set of samples collected from many episodes.

6 Experimental Study

This section evaluates the effectiveness of the proposed algorithms. The effectiveness of SOTD algorithm is illustrated via comparison to LSTD, which is also a batch TD algorithm. A comparison study of O^2TD is conducted with GTD2 and ETD as three off-policy convergent TD algorithms with linear computational cost per step.

6.1 Experimental Study of SOTD

The effectiveness of the proposed SOTD algorithm is shown by comparing the performance on the 400-state Random MDP domain [Dann *et al.*, 2014] with LSTD [Bradtke and Barto, 1996; Boyan, 1999] algorithm, which is one of the most sample-efficient algorithms to the best of our knowledge. Two widely used measurements in TD learning, Mean-Squares Projected Bellman Error (**MSPBE**) [Sutton *et al.*, 2009; Dann *et al.*, 2014] and Mean-Squares Error (**MSE**) are used as the error measurements.

This domain is a randomly generated MDP with 400 states and 10 actions [Dann *et al.*, 2014]. The transition probabilities are defined as $P(s'|s, a) \propto p_{ss'}^a + 10^{-5}$, where $p_{ss'}^a \sim U[0, 1]$. The behavior policy π_b , the target policy π as well as the start distribution are sampled in a similar manner. Each state is represented by a 201-dimensional feature vector, where 200 of the features were sampled from a uniform distribution, and the last feature was a constant one, the discount factor is set to $\gamma = 0.95$. The number of features $d = 200$, and we compare the performance of LSTD and SOTD with different numbers of training samples n , as shown in Figure 2. As Figure 2 shows, with relatively small sample size n , SOTD tends to be even more sample-efficient than the LSTD algorithm.

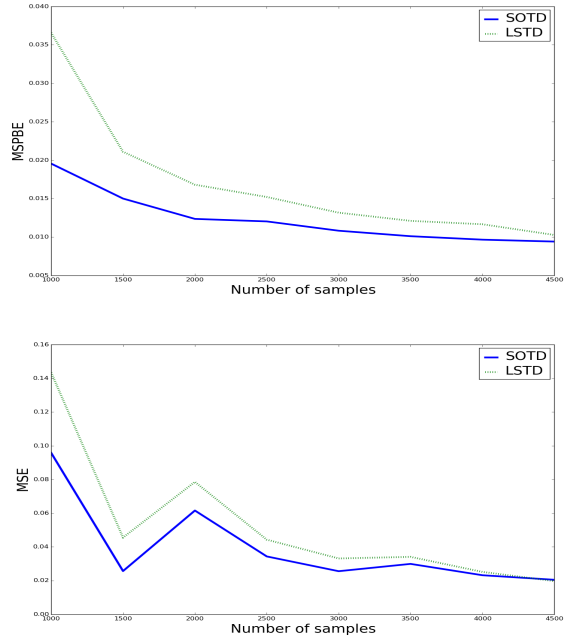


Figure 2: Comparison between SOTD and LSTD on 400-State Random MDP Domain

6.2 Experimental Study of O^2TD

This section compares the previous GTD2, ETD method with the O^2TD method using various domains with regard to their value function approximation performances. It should be mentioned that since the major focus of this paper is value function approximation and thus comparisons on control learning performance are not reported in this paper. We use α_E , α_O , and α_G to denote the stepsizes for ETD, O^2TD , and GTD2, respectively. Root Mean-Squares Projected Bellman Error (**RMSPBE**) and Root Mean-Squares Error (**RMSE**) are used for better visualization.

6.2.1 Baird Domain

The Baird example [Baird, 1995] is a well-known example to test the performance of off-policy convergent algorithms. Constant stepsize $\alpha_O = 0.006$, $\alpha_G = 0.005$, which are chosen via comparison studies as in [Dann *et al.*, 2014]. The Monte-Carlo estimation of true value function V is conducted as in [Dann *et al.*, 2014]. Figure 3 shows the RMSPBE curve and RMSE curve of GTD2, O^2TD of 5000 steps averaged over 20 runs. As can be seen from Figure 3, although the variance of O^2TD is larger than GTD2's, O^2TD has a significant improvement over the GTD2 algorithm wherein the RMSPBE, the RMSE and the variance are all substantially reduced. The low variance of the GTD2 learning curve can be explained by the advantage of stochastic gradient against stochastic approximation method, as

explained in Liu *et al.* [2015].

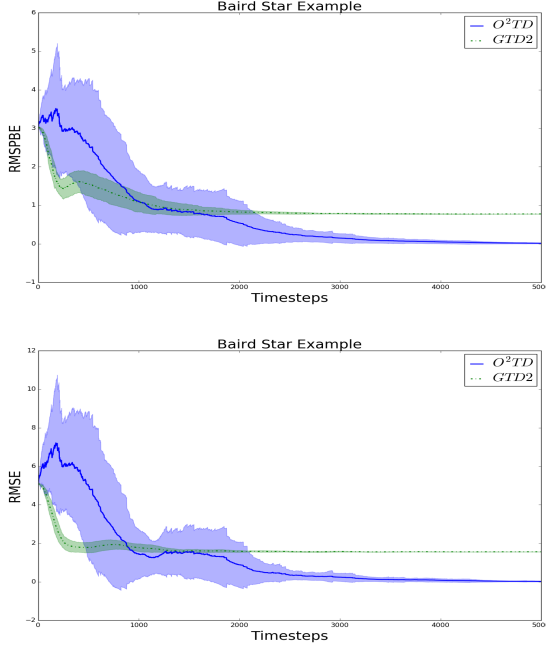


Figure 3: Baird Domain

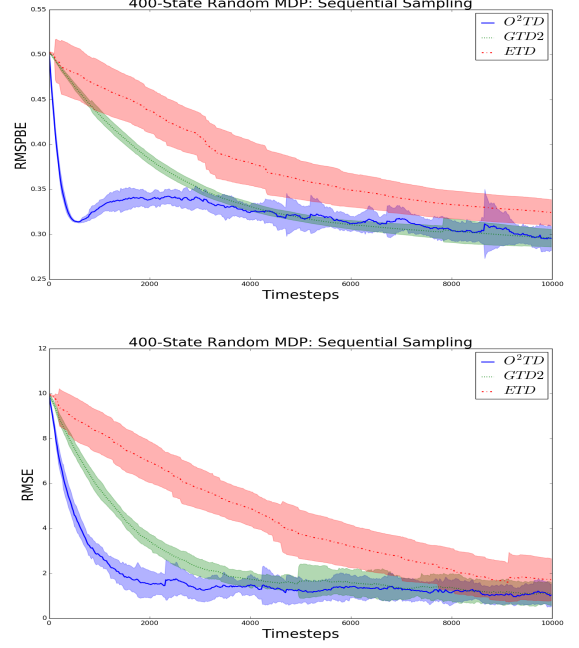


Figure 4: Random MDP with Sequential Sampling

6.2.2 400-State Random MDP

The randomly generated MDP with 400 states and 10 actions used in Section 6.1 is adopted as the second task. For sequential sampling (Figure 4), constant stepsize $\alpha_E = 3 * 10^{-6}$, $\alpha_O = 0.0007$, $\alpha_G = 0.002$. For random sampling (Figure 5), constant stepsize $\alpha_E = 2 * 10^{-6}$, $\alpha_O = 0.0006$, $\alpha_G = 0.0009$. The Monte-Carlo estimation of true value function V is conducted as in [Dann *et al.*, 2014]. ETD tends to diverge easily with large stepsizes on this domain, so α_E is set to be very small. As Figure 4 and Figure 5 show, O^2TD performs overall the best on this domain, although the variance is relatively larger than GTD2's.

6.2.3 Mountain Car

This section uses the mountain car example to evaluate the validity of the proposal algorithm. The mountain car MDP is an optimal control problem with a continuous two-dimensional state space. The step discontinuity in the value function makes learning difficult. The Fourier basis [Konidaris *et al.*, 2011] is used, which is a kind of fixed basis set. An empirically good policy π was obtained first, then we ran this policy π to collect trajectories that comprise the dataset. On-policy policy evaluation of π is then conducted using the collected samples. For sequential sampling, constant stepsize $\alpha_E = 0.001$, $\alpha_O = 0.1$, $\alpha_G = 0.2$.

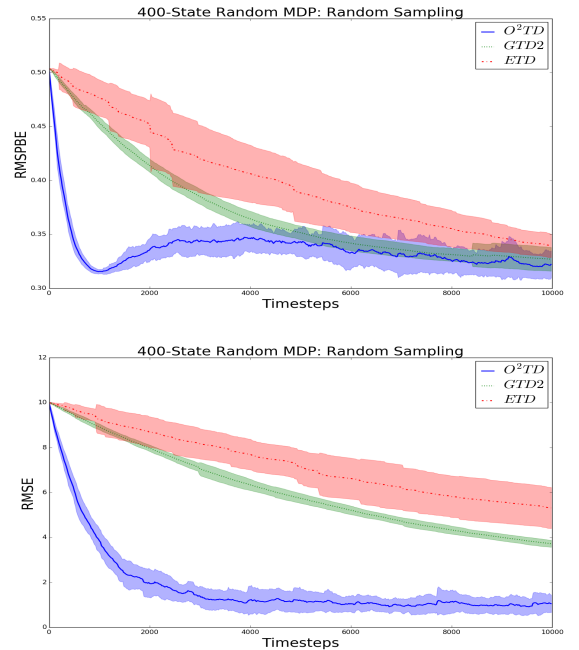


Figure 5: Random MDP with Random Sampling

For random sampling, constant stepsize $\alpha_E = 0.0002$, $\alpha_O = 0.05$, $\alpha_G = 0.06$. The Monte-Carlo estimation of V is estimated via 100 runs. As Figure 6 and Figure 7 show, GTD2 appears to perform the worst on this domain, and O^2TD tends to converge faster than ETD.

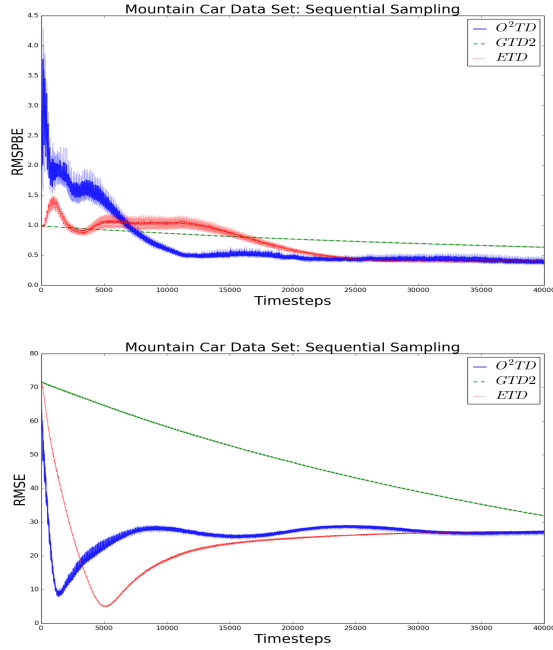


Figure 6: Mountain car with Sequential Sampling

7 Conclusion

This paper proposes an interesting question:

- How to improve the approximation quality of the true value function V ?

To this end, several algorithms are proposed that can apply to different scenarios. Empirical experimental studies solidify the effectiveness of the proposed algorithm with different learning settings.

The major contribution is *not* to propose another new TD algorithm with linear computational complexity per step, but to make an attempt to explore the optimal prediction of the value function in model-free policy evaluation. There are numerous promising future work potentials along this direction of research. One possible future research is to explore the relation between the near optimal projection matrix with eligibility traces and if the combination can improve the value function prediction performance in integration.

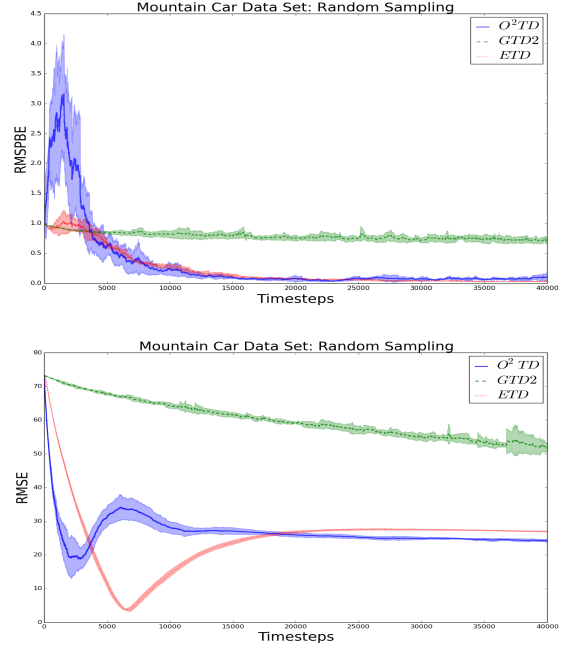


Figure 7: Mountain car with Random Sampling

Another interesting direction is that the current computationally tractable criteria of computing X^* are based on Proposition 1 and the power series expansion of $(L^\pi)^{-1}$, it would be very intriguing to explore if there exist other computationally tractable criteria.

References

- L. C. Baird. Residual algorithms: Reinforcement learning with function approximation. In *International Conference on Machine Learning*, pages 30–37, 1995.
- D. Bertsekas and J. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, Massachusetts, 1996.
- V. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge University Press, 2008.
- J. A. Boyan. Least-squares temporal difference learning. In *Proceedings of the 16th International Conference on Machine Learning*, pages 49–56. Morgan Kaufmann, San Francisco, CA, 1999.
- S. J. Bradtke and A. G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- C. Dann, G. Neumann, and J. Peters. Policy evaluation with temporal differences: A survey and comparison. *Journal of Machine Learning Research*, 15:809–883, 2014.

- G. J. Gordon. Stable fitted reinforcement learning. *Advances in neural information processing systems*, pages 1052–1058, 1996.
- A. Hallak, A. Tamar, R. Munos, and S. Mannor. Generalized emphatic temporal difference learning: Bias-variance analysis. *arXiv preprint arXiv:1509.05172*, 2015.
- G. Konidaris, S. Osentoski, and P. S. Thomas. Value function approximation in reinforcement learning using the fourier basis. In *Proceedings of the Twenty-Fifth Conference on Artificial Intelligence*, 2011.
- B. Liu, J. Liu, M. Ghavamzadeh, S. Mahadevan, and M. Petrik. Finite-sample analysis of proximal gradient td algorithms. In *Conference on Uncertainty in Artificial Intelligence*, 2015.
- B. Liu, L. Zhang, and J. Liu. Dantzig selector with an approximately optimal denoising matrix and its application in sparse reinforcement learning. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, pages 487–496, 2016.
- Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. 2004.
- B. Scherrer. Should one compute the temporal difference fix point or minimize the bellman residual? the unified oblique projection view. In *Proceedings of 27th International Conference on Machine Learning*, pages 52–68, 2010.
- S. P. Singh, T. Jaakkola, and M. I. Jordan. Reinforcement learning with soft state aggregation. *Advances in neural information processing systems*, pages 361–368, 1995.
- R. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- R. Sutton, C. Szepesvári, and H. Maei. A convergent $\mathcal{O}(n)$ algorithm for off-policy temporal-difference learning with linear function approximation. In *Neural Information Processing Systems*, pages 1609–1616, 2008.
- R. Sutton, H. Maei, D. Precup, S. Bhatnagar, D. Silver, C. Szepesvári, and E. Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *International Conference on Machine Learning*, pages 993–1000, 2009.
- R. S. Sutton, A. R. Mahmood, and M. White. An emphatic approach to the problem of off-policy temporal-difference learning. *The Journal of Machine Learning Research*, 17:1–29, 2015.
- J. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42:674–690, 1997.
- H. Yu. Convergence of least squares temporal difference methods under general conditions. In *Proceedings of the 27th International Conference on Machine Learning*, pages 1207–1214, 2010.
- H. Yu. Weak convergence properties of constrained emphatic temporal-difference learning with constant and slowly diminishing stepsize, 2015.

Appendix

Details of Eq. (24)

To obtain Eq. (24), we first introduce the following Lemmas to compute the singular value of rank-1 matrices.

Lemma 3. *A rank-1 real square matrix $G = pq^\top$ where p, q are vectors of the same length, the eigenvalues of G are*

$$\lambda(G) = \{p^\top q, 0, 0, 0, \dots\}, \quad (32)$$

i.e., G has only one nonzero eigenvalue $p^\top q$, and all other eigenvalues are 0, and thus we also have

$$\text{Tr}(G) = p^\top q, \quad (33)$$

where $\text{Tr}(\cdot)$ is the trace of a matrix.

Based on Lemma 3, we introduce Lemma 4.

Lemma 4. *A rank-1 real matrix (not necessarily to be square) $M = uv^\top$ has only one nonzero singular value $\sigma_{\max}(M) = \|u\|_2 \cdot \|v\|_2$, where $\|\cdot\|_2$ is the ℓ_2 -norm of a vector, and the Frobenius norm and the trace norm of M are identical, i.e.,*

$$\|M\|_* = \|M\|_F = \sigma_{\max}(M) = \|u\|_2 \cdot \|v\|_2 \quad (34)$$

Proof. We use M^H to represent the conjugate transpose of the M matrix, and $\lambda(\cdot)$ to represent the eigenvalues of a square matrix, and $\lambda(\cdot)$ to represent the nonzero eigenvalue of a matrix. Then we have

$$\begin{aligned} \lambda(M^H M) &= \lambda(vu^\top uv^\top) \\ &= (u^\top u) \lambda(vv^\top) \end{aligned} \quad (35)$$

From Lemma 3, we know that $\lambda(vv^\top)$ are $\{v^\top v, 0, 0, \dots\}$, and thus M has only one nonzero singular value $\sigma_{\max}(M)$, which is

$$\begin{aligned} \sigma_{\max}(M) &= \sqrt{\lambda(M^H M)} \\ &= \sqrt{\lambda(vu^\top uv^\top)} \\ &= \sqrt{(u^\top u) \lambda(vv^\top)} \\ &= \sqrt{(u^\top u)(v^\top v)} \\ &= \|u\|_2 \cdot \|v\|_2, \end{aligned}$$

and all other singular values of M are 0. Thus $\|M\|_* = \|M\|_F = \|u\|_2 \cdot \|v\|_2$, which completes the proof. \square

Based on Lemma 4, we now show the derivation of Eq. (24). To tackle the following trace norm minimization formulation,

$$\omega_i = \arg \min_{\omega} \|\omega \rho_i \phi_i \Delta \phi_i^\top - \phi_i \phi_i^\top\|_*, \quad (36)$$

we need to utilize the structure of the rank-1 matrices. We have

$$\omega \rho_i \phi_i \Delta \phi_i^\top - \phi_i \phi_i^\top = \phi_i (\omega \rho_i \Delta \phi_i - \phi_i)^\top, \quad (37)$$

we denote $q_i(\omega) := (\omega \rho_i \Delta \phi_i - \phi_i)^\top$, and thus we have

$$\begin{aligned} \omega_i &= \arg \min_{\omega} \|\phi_i q_i^\top(\omega)\|_* \\ &= \arg \min_{\omega} \|\phi_i\|_2 \cdot \|q_i(\omega)\|_2 \\ &= \arg \min_{\omega} \|q_i(\omega)\|_2 \end{aligned} \quad (38)$$

The second equality comes based on Eq. (34), and the third equality is based on the fact that $\|\phi_i\|_2$ does not depend on ω . This is equivalent to the following,

$$\omega_i = \arg \min_{\omega} \|\omega \rho_i \Delta \phi_i - \phi_i\|_2^2 \quad (39)$$

On the other hand, if we use $\|\cdot\|_F^2$ instead of trace norm minimization as in Eq. (36), we have

$$\omega_i = \arg \min_{\omega} \|\phi_i q_i(\omega)\|_F^2, \quad (40)$$

And since

$$\begin{aligned} \|\phi_i q_i^\top(\omega)\|_F^2 &= \text{Tr}(q_i(\omega) \phi_i^\top \phi_i q_i^\top(\omega)) \\ &= (\phi_i^\top \phi_i) \text{Tr}(q_i(\omega) q_i^\top(\omega)) \\ &= (\phi_i^\top \phi_i) (q_i^\top(\omega) q_i(\omega)) \\ &= \|\phi_i\|_2^2 \|q_i(\omega)\|_2^2. \end{aligned} \quad (41)$$

The first equality comes from that for a matrix M , there is

$$\|M\|_F^2 = \text{Tr}(M^H M). \quad (42)$$

The third equality comes from Eq. (39). Then we can see that problem (40) is also equivalent to Eq. (39), as verified by Lemma 4.

By taking the gradient of the right hand-side of Eq. (39), we will have Eq. (24) as the final result.